# Clustering
## (unsupervised learning)

JJ Valletta

March 4, 2015

# Overview

- What is clustering?

- Major types of clustering methods

- $k$-means clustering

- Agglomerative hierarchical clustering

- Gaussian mixture models

- How do we determine the correct number of clusters?

# Overview

- ## What is clustering?

- Major types of clustering methods

- $k$-means clustering

- Agglomerative hierarchical clustering

- Gaussian mixture models

- How do we determine the correct number of clusters?

# Overview

- What is clustering?

- Major types of clustering methods

- $k$-means clustering

- Agglomerative hierarchical clustering

- Gaussian mixture models

- How do we determine the correct number of clusters?

# Overview

- What is clustering?

- Major types of clustering methods

- $k$-means clustering

- Agglomerative hierarchical clustering

- Gaussian mixture models

- How do we determine the correct number of clusters?

# Overview

- What is clustering?

- Major types of clustering methods

- $k$-means clustering

- Agglomerative hierarchical clustering

- Gaussian mixture models

- How do we determine the correct number of clusters?

# Overview

- What is clustering?

- Major types of clustering methods

- $k$-means clustering

- Agglomerative hierarchical clustering

- Gaussian mixture models

- How do we determine the correct number of clusters?

# Overview

- What is clustering?

- Major types of clustering methods

- $k$-means clustering

- Agglomerative hierarchical clustering

- Gaussian mixture models

- How do we determine the correct number of clusters?

# What is clustering?

# What is clustering?

## Formal definition

Identifying homogeneous and well separated groups of data points (features) by some similarity measure

# What is clustering?

## Formal definition

Identifying homogeneous and well separated groups of data points (features) by some similarity measure

## Informal definition

The process of stereotyping your data
e.g *these* are round(ish) faces, *these* are short(ish) people

# What is clustering?

## Formal definition

Identifying homogeneous and well separated groups of data points (features) by some similarity measure

## Informal definition

The process of stereotyping your data
e.g *these* are round(ish) faces, *these* are short(ish) people

## But how many groups?

An unsolved problem. Issue lies in the subjectivity of the word **similar** and its mathematical definition
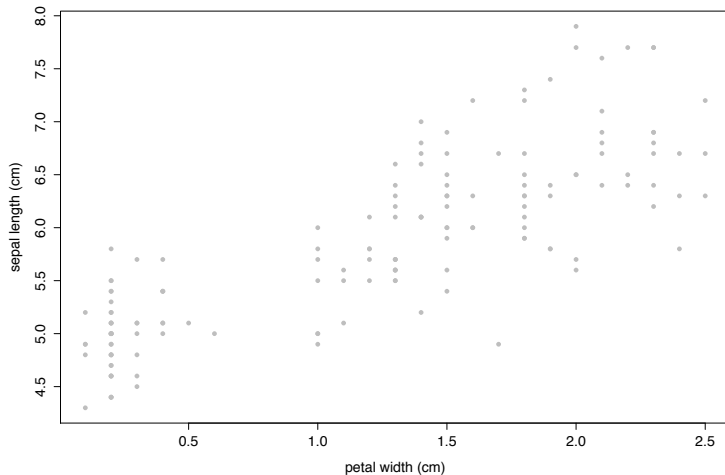
# Are they similar?

# Are they similar?

# What are we after?

- High intra-cluster similarity
- Low inter-cluster similarity
- Elucidate on how the data is structured (maybe identify outliers)

# Where is clustering used?

Biological systematics: finding organisms sharing similar attributes

Computer vision: segmenting a digital image for object recognition

Epidemiology: identifying geographical clusters of diseases

Gene expression: discovering co-regulated genes

Medical imaging: differentiating between tissues

Mathematical chemistry: grouping compounds by topological indices

**Clustering is particularly useful in applications where labelling the data is very time consuming/expensive**
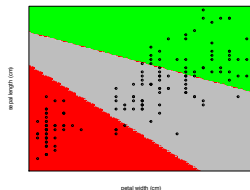
# Major types of clustering methods

Partitional: The data (feature) space is partitioned into $k$ regions

Hierarchical: Iteratively merging small clusters into larger ones (*agglomerative*) or breaking large clusters into smaller ones (*divisive*)

Distribution-based: Fit $k$ multivariate statistical distributions

# Major types of clustering methods

Partial: The data (feature) space is
partitioned into $k$ regions



Hierarchical: Iteratively merging small
clusters into larger ones
(*agglomerative*) or breaking
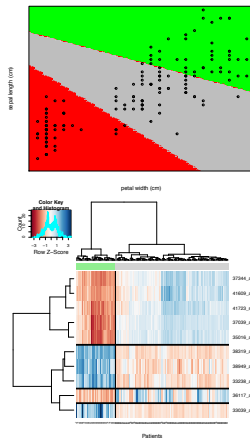large clusters into smaller
ones (*divisive*)

Distribution-based: Fit $k$ multivariate statistical
distributions

# Major types of clustering methods

Partial: The data (feature) space is partitioned into $k$ regions



Hierarchical: Iteratively merging small clusters into larger ones (*agglomerative*) or breaking large clusters into smaller ones (*divisive*)



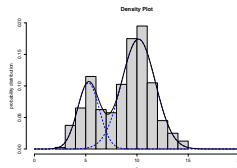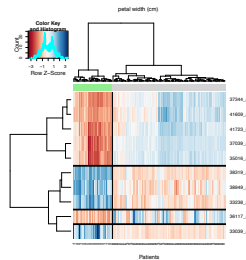Distribution-based: Fit $k$ multivariate statistical distributions

# Major types of clustering methods

Partitional: The data (feature) space is partitioned into $k$ regions



Hierarchical: Iteratively merging small clusters into larger ones (*agglomerative*) or breaking large clusters into smaller ones (*divisive*)



Distribution-based: Fit $k$ multivariate statistical distributions

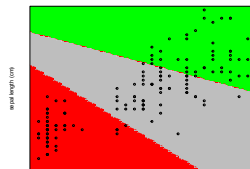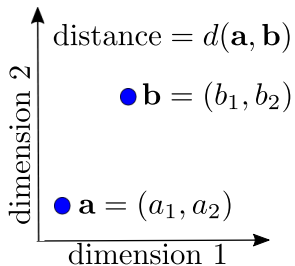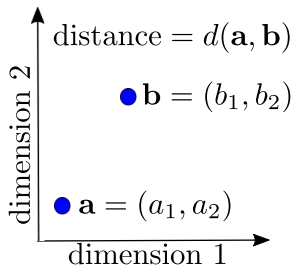# Similarity measures

What is the distance $d(\mathbf{a}, \mathbf{b})$ between $\mathbf{a}$ and $\mathbf{b}$?

distance $= d(\mathbf{a}, \mathbf{b})$

$\bullet \, \mathbf{b} = (b_1, b_2)$

$\bullet \, \mathbf{a} = (a_1, a_2)$

dimension 2

dimension 1

# Similarity measures

What is the distance $d(\mathbf{a}, \mathbf{b})$ between $\mathbf{a}$ and $\mathbf{b}$?



| example: $\mathbf{a}, \mathbf{b} \in \mathbb{R}^2$ | in general: $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ |
|---|---|
| • **Manhattan** $\|a_1 - b_1\| + \|a_2 - b_2\|$ | $\sum_{i=1}^{d} \|a_i - b_i\| = \|\mathbf{a} - \mathbf{b}\|_1$ |
| • Euclidean $\sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2}$ | $\left(\sum_{i=1}^{d} (a_i - b_i)^2\right)^{\frac{1}{2}} = \|\mathbf{a} - \mathbf{b}\|_2$ |
| • Minkowski (p-norm) $\sqrt[p]{\|a_1 - b_1\|^p + \|a_2 - b_2\|^p}$ | $\left(\sum_{i=1}^{d} \|a_i - b_i\|^p\right)^{\frac{1}{p}} = \|\mathbf{a} - \mathbf{b}\|_p$ |

## Similarity measures

What is the distance $d(\mathbf{a}, \mathbf{b})$ between $\mathbf{a}$ and $\mathbf{b}$?



distance $= d(\mathbf{a}, \mathbf{b})$

$\bullet \, \mathbf{b} = (b_1, b_2)$

$\bullet \, \mathbf{a} = (a_1, a_2)$

dimension 2

dimension 1

| example: $\mathbf{a}, \mathbf{b} \in \mathbb{R}^2$ | in general: $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ |
|---|---|
| • **Manhattan** $|a_1 - b_1| + |a_2 - b_2|$ | $\sum_{i=1}^{d} |a_i - b_i| = \|\mathbf{a} - \mathbf{b}\|_1$ |
| • **Euclidean** $\sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2}$ | $\left(\sum_{i=1}^{d}(a_i - b_i)^2\right)^{\frac{1}{2}} = \|\mathbf{a} - \mathbf{b}\|_2$ |
| • **Minkowski (p-norm)** $\sqrt[p]{|a_1 - b_1|^p + |a_2 - b_2|^p}$ | $\left(\sum_{i=1}^{d} |a_i - b_i|^p\right)^{\frac{1}{p}} = \|\mathbf{a} - \mathbf{b}\|_p$ |

## Similarity measures

What is the distance $d(\mathbf{a}, \mathbf{b})$ between $\mathbf{a}$ and $\mathbf{b}$?



| example: $\mathbf{a}, \mathbf{b} \in \mathbb{R}^2$ | in general: $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ |
|---|---|
| • **Manhattan** $\|a_1 - b_1\| + \|a_2 - b_2\|$ | $\sum_{i=1}^{d} \|a_i - b_i\| = \|\mathbf{a} - \mathbf{b}\|_1$ |
| • **Euclidean** $\sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2}$ | $\left(\sum_{i=1}^{d}(a_i - b_i)^2\right)^{\frac{1}{2}} = \|\mathbf{a} - \mathbf{b}\|_2$ |
| • **Minkowski (p-norm)** $\sqrt[p]{\|a_1 - b_1\|^p + \|a_2 - b_2\|^p}$ | $\left(\sum_{i=1}^{d} \|a_i - b_i\|^p\right)^{\frac{1}{p}} = \|\mathbf{a} - \mathbf{b}\|_p$ |

# Similarity measures

What is the distance $d(\mathbf{a},\ \mathbf{b})$ between $\mathbf{a}$ and $\mathbf{b}$?



distance $= d(\mathbf{a}, \mathbf{b})$

$\bullet\, \mathbf{b} = (b_1, b_2)$

$\bullet\, \mathbf{a} = (a_1, a_2)$

dimension 2

dimension 1

## example: $\mathbf{a},\ \mathbf{b} \in \mathbb{R}^2$

- **Canberra**
  $\frac{|a_1 - b_1|}{|a_1| + |b_1|} + \frac{|a_2 - b_2|}{|a_2| + |b_2|}$

- **Cosine similarity**
  $\frac{a_1 b_1 + a_2 b_2}{\sqrt{a_1^2 + a_2^2}\sqrt{b_1^2 + b_2^2}}$

- **Correlation distance**

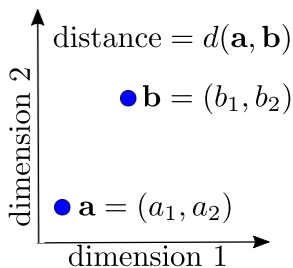## in general: $\mathbf{a},\ \mathbf{b} \in \mathbb{R}^d$

$\sum_{i=1}^{d} \frac{|a_i - b_i|}{|a_i| + |b_i|}$

$\frac{\sum_{i=1}^{d} a_i b_i}{\sqrt{\sum_{i=1}^{d}(a_i)^2}\sqrt{\sum_{i=1}^{d}(b_i)^2}} = \frac{\mathbf{a}.\mathbf{b}}{\|a\|_2\|b\|_2}$

Pearson ($\rho$), Spearman, Kendall ($\tau$)

# Similarity measures

What is the distance $d(\mathbf{a}, \mathbf{b})$ between $\mathbf{a}$ and $\mathbf{b}$?



| example: $\mathbf{a}, \mathbf{b} \in \mathbb{R}^2$ | in general: $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ |
|---|---|
| • **Canberra** $\frac{|a_1-b_1|}{|a_1|+|b_1|} + \frac{|a_2-b_2|}{|a_2|+|b_2|}$ | $\sum_{i=1}^{d} \frac{|a_i-b_i|}{|a_i|+|b_i|}$ |
| • Cosine similarity $\frac{a_1 b_1 + a_2 b_2}{\sqrt{a_1^2+a_2^2}\sqrt{b_1^2+b_2^2}}$ | $\frac{\sum_{i=1}^{d} a_i b_i}{\sqrt{\sum_{i=1}^{d}(a_i)^2}\sqrt{\sum_{i=1}^{d}(b_i)^2}} = \frac{\mathbf{a}.\mathbf{b}}{\|a\|_2\|b\|_2}$ |
| • Correlation distance | Pearson ($\rho$), Spearman, Kendall ($\tau$) |

# Similarity measures

What is the distance $d(\mathbf{a}, \ \mathbf{b})$ between $\mathbf{a}$ and $\mathbf{b}$?



## example: $\mathbf{a}, \ \mathbf{b} \in \mathbb{R}^2$

- **Canberra**
  $\frac{|a_1-b_1|}{|a_1|+|b_1|} + \frac{|a_2-b_2|}{|a_2|+|b_2|}$

- **Cosine similarity**
  $\frac{a_1 b_1 + a_2 b_2}{\sqrt{a_1^2 + a_2^2}\sqrt{b_1^2 + b_2^2}}$

- **Correlation distance**

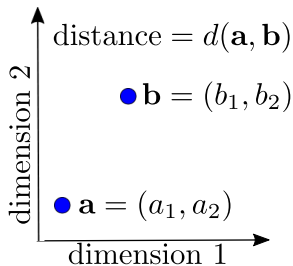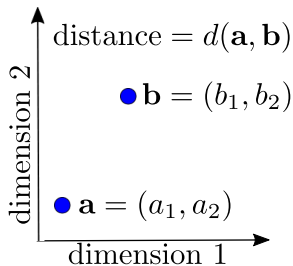## in general: $\mathbf{a}, \ \mathbf{b} \in \mathbb{R}^d$

$\sum_{i=1}^{d} \frac{|a_i - b_i|}{|a_i| + |b_i|}$

$\frac{\sum_{i=1}^{d} a_i b_i}{\sqrt{\sum_{i=1}^{d}(a_i)^2}\sqrt{\sum_{i=1}^{d}(b_i)^2}} = \frac{\mathbf{a}.\mathbf{b}}{\|a\|_2 \|b\|_2}$

Pearson ($\rho$), Spearman, Kendall ($\tau$)

# Similarity measures

What is the distance $d(\mathbf{a}, \mathbf{b})$ between $\mathbf{a}$ and $\mathbf{b}$?



| example: $\mathbf{a}, \mathbf{b} \in \mathbb{R}^2$ | in general: $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ |
|---|---|
| • **Canberra** $\frac{\lvert a_1 - b_1 \rvert}{\lvert a_1 \rvert + \lvert b_1 \rvert} + \frac{\lvert a_2 - b_2 \rvert}{\lvert a_2 \rvert + \lvert b_2 \rvert}$ | $\sum_{i=1}^{d} \frac{\lvert a_i - b_i \rvert}{\lvert a_i \rvert + \lvert b_i \rvert}$ |
| • **Cosine similarity** $\frac{a_1 b_1 + a_2 b_2}{\sqrt{a_1^2 + a_2^2} \sqrt{b_1^2 + b_2^2}}$ | $\frac{\sum_{i=1}^{d} a_i b_i}{\sqrt{\sum_{i=1}^{d} (a_i)^2} \sqrt{\sum_{i=1}^{d} (b_i)^2}} = \frac{\mathbf{a}.\mathbf{b}}{\lVert a \rVert_2 \lVert b \rVert_2}$ |
| • **Correlation distance** | Pearson ($\rho$), Spearman, Kendall ($\tau$) |

# $k$-means clustering

```
fit <- kmeans(x, centers)
# x - numeric matrix of data
# centers - no. of clusters k
```

1. Select $k$ centroids at random

2. Calculate distance between centroids and each data point

3. Assign each data point to the closest centroid

4. Compute new centroids; the average of all data points in that cluster

5. Repeat steps 2 to 4 until data points remain in the same cluster or some maximum number of iterations reached

# $k$-means clustering

```
fit <- kmeans(x, centers)
# x - numeric matrix of data
# centers - no. of clusters k
```

1. Select $k$ centroids at random

2. Calculate distance between centroids and each data point

3. Assign each data point to the closest centroid

4. Compute new centroids; the average of all data points in that cluster

5. Repeat steps 2 to 4 until data points remain in the same cluster or some maximum number of iterations reached

# $k$-means clustering

```
fit <- kmeans(x, centers)
# x - numeric matrix of data
# centers - no. of clusters k
```

**①** Select $k$ centroids at random

**②** Calculate distance between centroids and each data point

**③** Assign each data point to the closest centroid

**④** Compute new centroids; the average of all data points in that cluster

**⑤** Repeat steps 2 to 4 until data points remain in the same cluster or some maximum number of iterations reached

# $k$-means clustering

```
fit <- kmeans(x, centers)
# x - numeric matrix of data
# centers - no. of clusters k
```

1. Select $k$ centroids at random

2. Calculate distance between centroids and each data point

3. Assign each data point to the closest centroid

4. Compute new centroids; the average of all data points in that cluster

5. Repeat steps 2 to 4 until data points remain in the same cluster or some maximum number of iterations reached

# $k$-means clustering

```
fit <- kmeans(x, centers)
# x - numeric matrix of data
# centers - no. of clusters k
```

1. Select $k$ centroids at random

2. Calculate distance between centroids and each data point

3. Assign each data point to the closest centroid

4. Compute new centroids; the average of all data points in that cluster

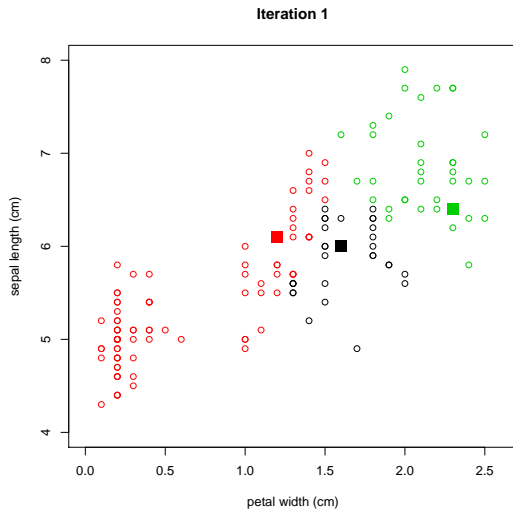5. Repeat steps 2 to 4 until data points remain in the same cluster or some maximum number of iterations reached

# $k$-means clustering

```
fit <- kmeans(x, centers)
# x - numeric matrix of data
# centers - no. of clusters k
```
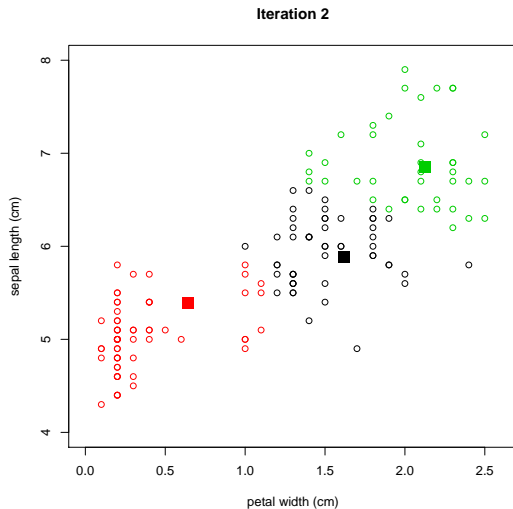
1. Select $k$ centroids at random

2. Calculate distance between centroids and each data point

3. Assign each data point to the closest centroid

4. Compute new centroids; the average of all data points in that cluster

5. Repeat steps 2 to 4 until data points remain in the same cluster or some maximum number of iterations reached

# $k$-means clustering

```
fit <- kmeans(x, centers)
# x - numeric matrix of data
# centers - no. of clusters k
```

1. Select $k$ centroids at random

2. Calculate distance between centroids and each data point

3. Assign each data point to the closest centroid

4. Compute new centroids; the average of all data points in that cluster

5. Repeat steps 2 to 4 until data points remain in the same cluster or some maximum number of iterations reached

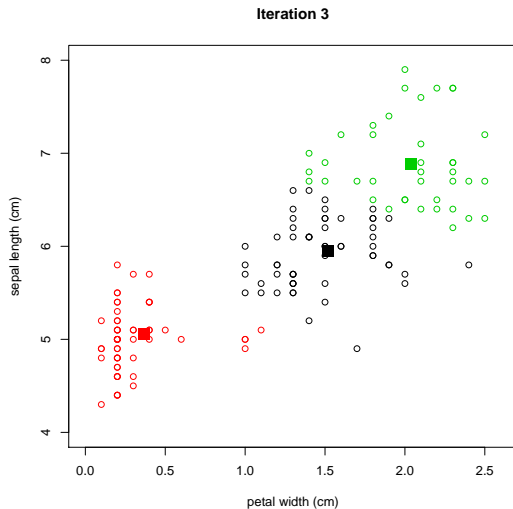**Note**: $k$-means clustering should *only* be used with continuous data
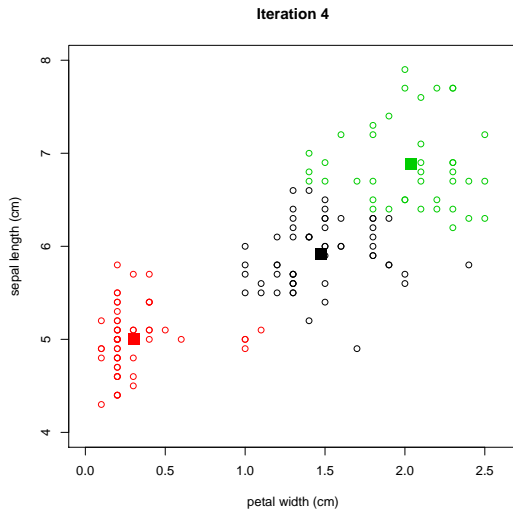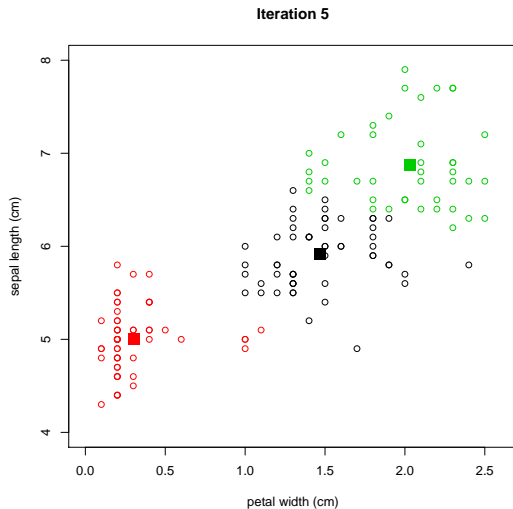
# $k$-means clustering

# $k$-means clustering
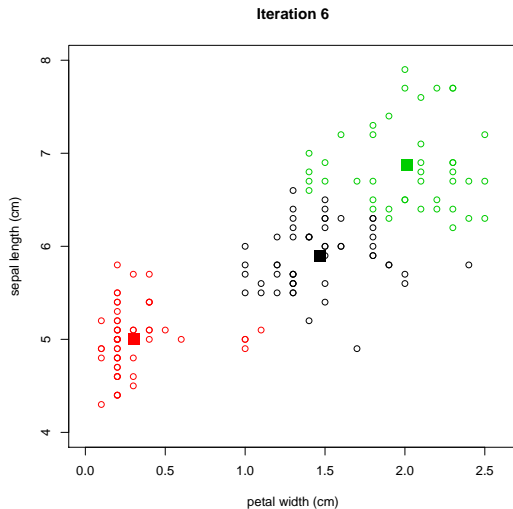
# $k$-means clustering
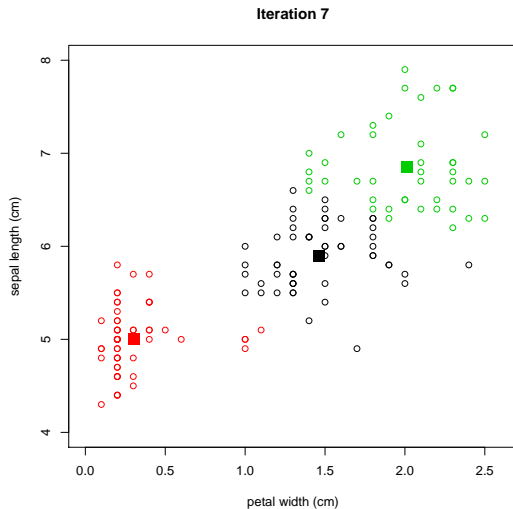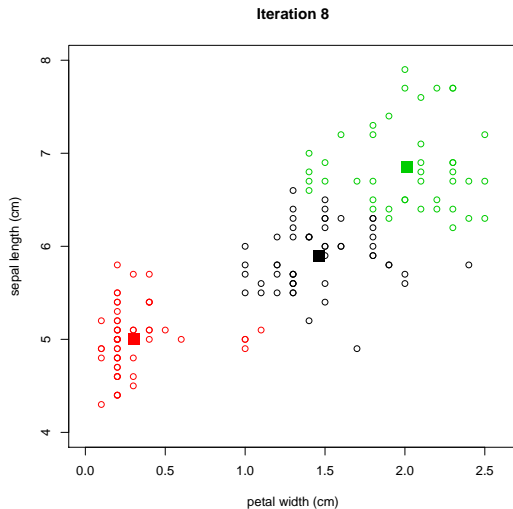
# $k$-means clustering

# $k$-means clustering

# $k$-means clustering

# $k$-means clustering

# $k$-means clustering

# $k$-means clustering

## Pros

- Simple and intuitive
- Computationally inexpensive/fast

## Cons

- What is $k$?
- Only applicable to continuous data where a mean is defined
- No guarantee of a global optimum solution

# Agglomerative hierarchical clustering

```
d <- dist(as.matrix(data), method)
# data - data frame
# method - distance method e.g "euclidean" or "manhattan"
fit <- hclust(d, method)
# method - linkage function e.g "complete" or "single"
```

1. Assign each data point as its own cluster

2. Compute distance between each cluster

3. Merge the closest pair into a single cluster

4. Repeat 2 to 3 until you're left with one cluster

# Agglomerative hierarchical clustering

```
d <- dist(as.matrix(data), method)
# data - data frame
# method - distance method e.g "euclidean" or "manhattan"
fit <- hclust(d, method)
# method - linkage function e.g "complete" or "single"
```

1. Assign each data point as its own cluster

2. Compute distance between each cluster

3. Merge the closest pair into a single cluster

4. Repeat 2 to 3 until you're left with one cluster

# Agglomerative hierarchical clustering

```
d <- dist(as.matrix(data), method)
# data - data frame
# method - distance method e.g "euclidean" or "manhattan"
fit <- hclust(d, method)
# method - linkage function e.g "complete" or "single"
```

1. Assign each data point as its own cluster

2. Compute distance between each cluster

3. Merge the closest pair into a single cluster

4. Repeat 2 to 3 until you're left with one cluster

# Agglomerative hierarchical clustering

```
d <- dist(as.matrix(data), method)
# data - data frame
# method - distance method e.g "euclidean" or "manhattan"
fit <- hclust(d, method)
# method - linkage function e.g "complete" or "single"
```

1. Assign each data point as its own cluster

2. Compute distance between each cluster

3. Merge the closest pair into a single cluster

4. Repeat 2 to 3 until you're left with one cluster

# Agglomerative hierarchical clustering

```
d <- dist(as.matrix(data), method)
# data - data frame
# method - distance method e.g "euclidean" or "manhattan"
fit <- hclust(d, method)
# method - linkage function e.g "complete" or "single"
```

1. Assign each data point as its own cluster

2. Compute distance between each cluster

3. Merge the closest pair into a single cluster

4. Repeat 2 to 3 until you're left with one cluster

# Agglomerative hierarchical clustering

```
d <- dist(as.matrix(data), method)
# data - data frame
# method - distance method e.g "euclidean" or "manhattan"
fit <- hclust(d, method)
# method - linkage function e.g "complete" or "single"
```

1. Assign each data point as its own cluster

2. Compute distance between each cluster

3. Merge the closest pair into a single cluster

4. Repeat 2 to 3 until you're left with one cluster

**Note**: Step 3 is *key*, the distance method and linkage function dictate the final result

# Hierarchical clustering: Link method

How do we calculate the inter-cluster distance? The *linkage function*

Centroid: mean of data points (same as in
$k$-means)

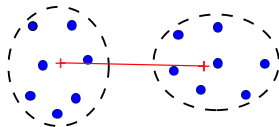Single: distance between closest pair of
points

Complete: distance between furthest pair of
points

Average: mean pairwise distance between
all points

# Hierarchical clustering: Link method

How do we calculate the inter-cluster distance? The *linkage function*

Centroid: mean of data points (same as in $k$-means)



Single: distance between closest pair of points

Complete: distance between furthest pair of points

Average: mean pairwise distance between all points

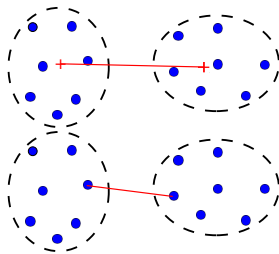# Hierarchical clustering: Link method

How do we calculate the inter-cluster distance? The *linkage function*

Centroid: mean of data points (same as in $k$-means)

Single: distance between closest pair of points

Complete: distance between furthest pair of points

Average: mean pairwise distance between all points

# Hierarchical clustering: Link method

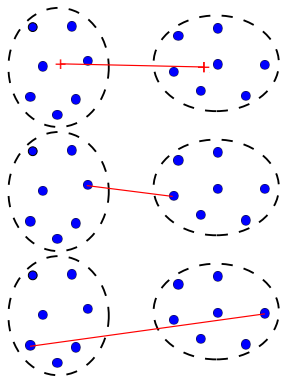How do we calculate the inter-cluster distance? The *linkage function*

Centroid: mean of data points (same as in $k$-means)

Single: distance between closest pair of points

Complete: distance between furthest pair of points

Average: mean pairwise distance between all points
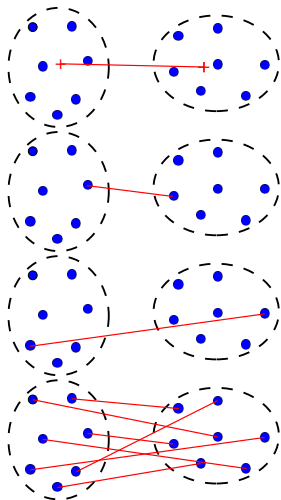
# Hierarchical clustering: Link method

How do we calculate the inter-cluster distance? The *linkage function*
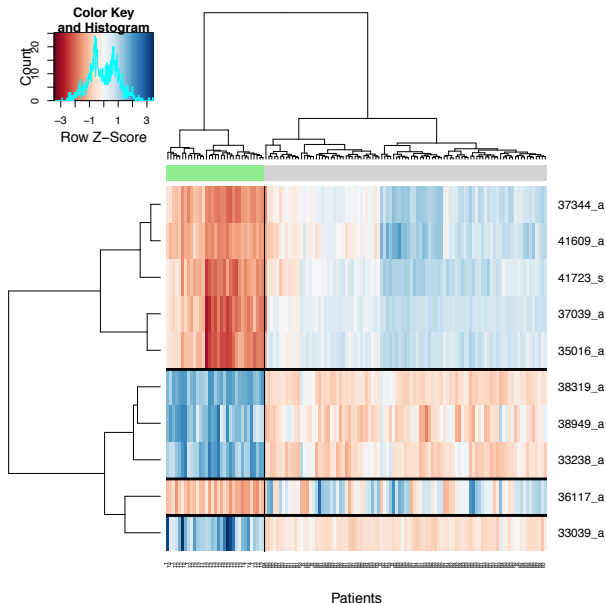


Centroid: mean of data points (same as in $k$-means)

Single: distance between closest pair of points

Complete: distance between furthest pair of points

Average: mean pairwise distance between all points

# Hierarchical clustering in gene expression studies

# Hierarchical clustering

## Pros

- No need to specify $k$
- Results can be visualised nicely irrespective of number of dimensions

## Cons

- Can be computationally expensive
- Interpretation is subjective. Where should we draw the line (to separate clusters)?
- Choice of distance method and linkage function can significantly change the result

# Gaussian mixture models

```
library(mclust)
fit <- Mclust(data, G)
# data - data frame
# G - no. of Gaussians
```
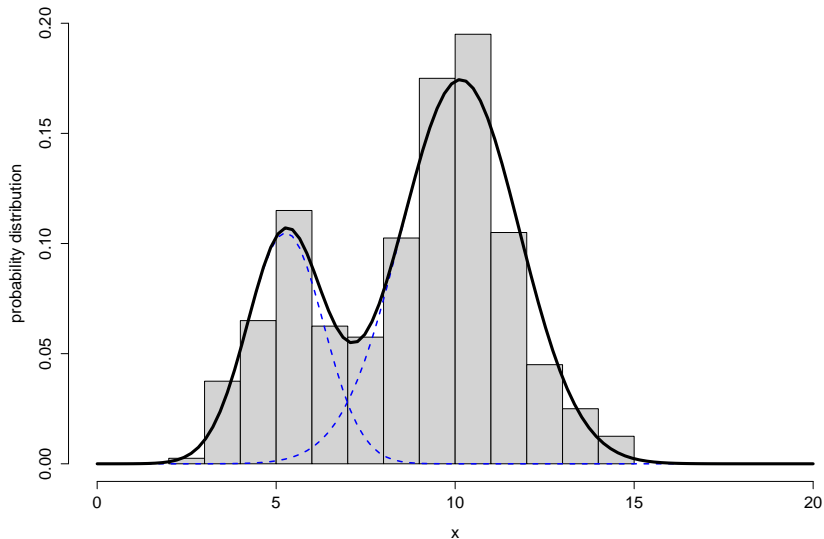
**1** Fit $k$ multivariate Gaussian distributions

The Expectation-Maximisation (EM) algorithm is used to estimate the parameters $\pi_i$ (mixing coefficients), $\mu_i$ and $\sigma_i$

$$p(x) = \sum_{i=1}^{k} \pi_i \mathcal{N}(x|\mu_i, \Sigma_i) \text{ and } \sum_{i=1}^{k} \pi_i = 1$$

Can be seen as a "soft" version of $k$-means because *every* point is part of *every* cluster but with varying levels of membership
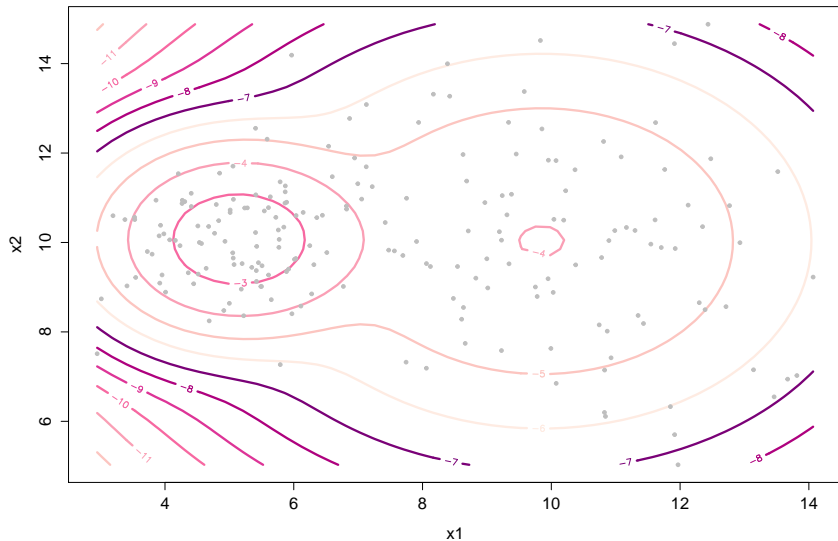
# Gaussian mixture models



**Density Plot**

# Gaussian mixture models



**log Density Contour Plot**

# Gaussian mixture models

## Pros

- Intuitive interpretation
- Computationally inexpensive

## Cons

- What is $k$?
- Strong assumption on the data (normality)
- No guarantee of a global optimum solution

## How do we determine the correct number of clusters?

**Short answer**: you can't

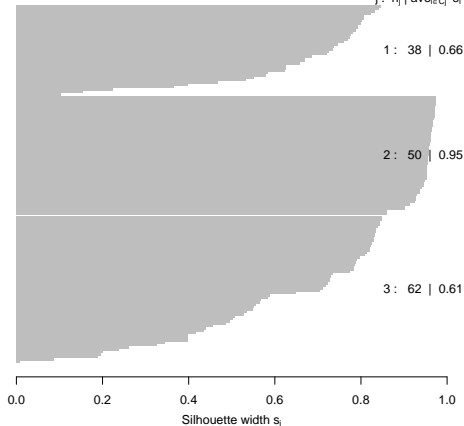Because data is unlabelled the correct number of $k$ is ambiguous

However we can plot some indices as a function of $k$ to help us evaluate cluster validity:

- Within cluster sum of square distances

- Akaike Information Criterion (AIC) or Bayesian Information Criterion (BIC) when using distribution-based methods

- Silhouette plot $-1 \geq s(i) \leq 1$ where:
  $s(i) = 1$, $i$th datum is appropriately clustered (good)
  $s(i) = 0$, $i$th datum is borderline between two clusters (meh.)
  $s(i) = -1$, $i$th datum should be in neighbouring cluster (bad)
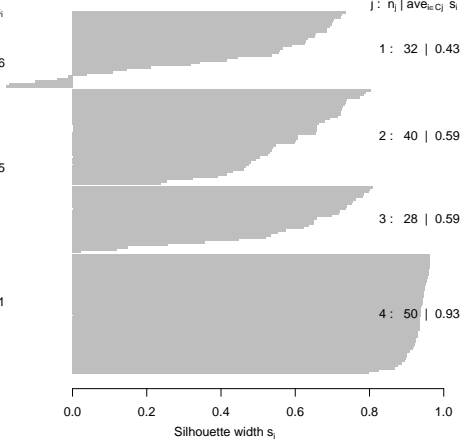
# How do we determine the correct number of clusters?



**Silhouette plot k=3**

n = 150

3 clusters $C_j$
j : $n_j$ | ave$_{i \in C_j}$ $s_i$

1 : 38 | 0.66

2 : 50 | 0.95

3 : 62 | 0.61

0.0    0.2    0.4    0.6    0.8    1.0

Silhouette width $s_i$

Average silhouette width : 0.74

**Silhouette plot k=4**

n = 150

4 clusters $C_j$
j : $n_j$ | ave$_{i \in C_j}$ $s_i$

1 : 32 | 0.43

2 : 40 | 0.59

3 : 28 | 0.59

4 : 50 | 0.93

0.0    0.2    0.4    0.6    0.8    1.0

Silhouette width $s_i$

Average silhouette width : 0.67

# How do we determine the correct number of clusters?

The `NbClust` package provides 30 different cluster validity metrics. A
majority vote can be taken to deduce the appropriate number of clusters

```
library(NbClust)
NbClust(data, distance, method, min.nc, max.nc, index)
# data - data frame
# distance - similarity measure e.g "euclidean"
# method - clustering algorithm e.g "kmeans"
# min.nc - min number of clusters to consider
# max.nc - max number of clusters to consider
# index - which indices to compute, "all" computes all of them
```

**Note**: These indices can *only* give us a ballpark range for the correct
number of clusters